

1.

2.

$$err(\mathbf{w}) = \begin{cases} 0, & -y\mathbf{w}^T \mathbf{x} \leq 0 (\text{correct classified}) \\ -y\mathbf{w}^T \mathbf{x}, & -y\mathbf{w}^T \mathbf{x} > 0 (\text{incorrect classified}) \end{cases}$$

Here we only consider the case that some points are classified wrong, meaning that we only discuss cases that make $err(\mathbf{w}) = -y\mathbf{w}^T \mathbf{x}$.

$$\frac{\partial err(\mathbf{w})}{\partial w_i} = -yx_i$$

While w_i stands for the i^{th} element in vector \mathbf{w} , x_i stands for the i^{th} element in vector \mathbf{x} .

So, the gradient of $-y\mathbf{w}^T \mathbf{x}$ is:

$$\nabla(-y\mathbf{w}^T \mathbf{x}) = -y\mathbf{x}$$

Therefore, for a point that is classified wrong, according to the update rules of SGD,

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla(-y\mathbf{w}^T \mathbf{x}) = \mathbf{w}_t + \eta y\mathbf{x}, \eta \text{ is the step size}$$

While the update rule of PLA looks like this when a point is classified wrong:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + y\mathbf{x}$$

So when $\eta = 1$, $err(\mathbf{w}) = \max(0, -y\mathbf{w}^T \mathbf{x})$, using SGD will result in PLA.

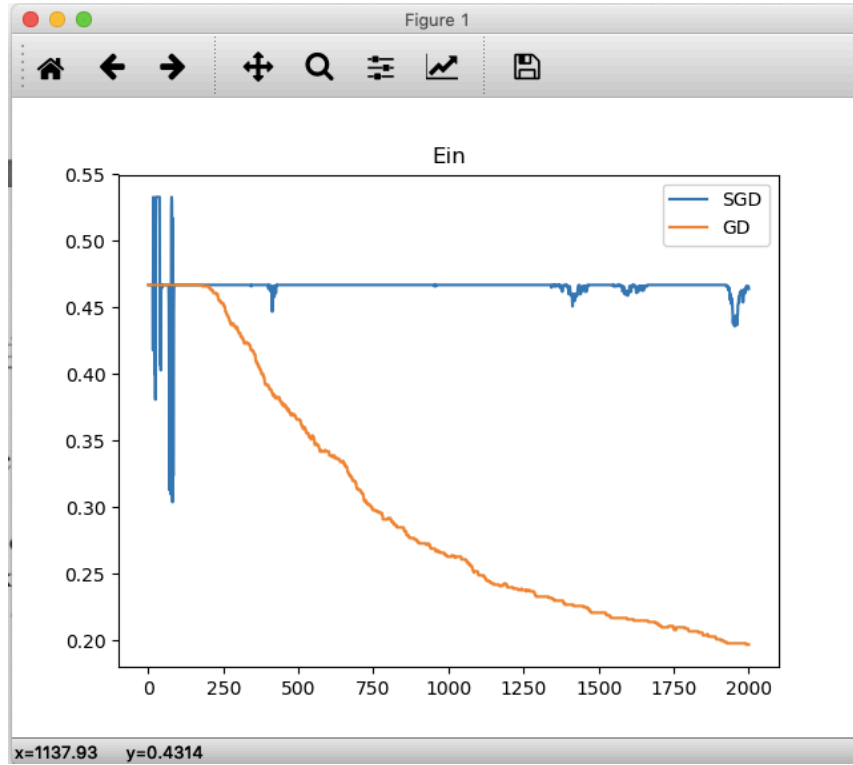
3.

$$E_{in} = \frac{1}{N} \sum_{n=1}^N \left(\ln \left(\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n) \right) - \mathbf{w}_{y_n}^T \mathbf{x}_n \right)$$

$$\frac{\partial E_{in}}{\partial w_i} = \frac{1}{N} \sum_{n=1}^N \left(\frac{\partial \left(\ln \left(\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n) \right) \right)}{\partial w_i} - \frac{\partial (\mathbf{w}_{y_n}^T \mathbf{x}_n)}{\partial w_i} \right)$$

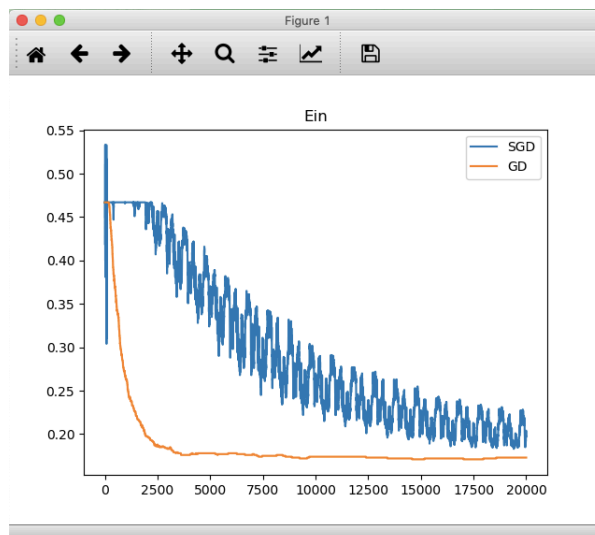
$$\begin{aligned}
&= \frac{1}{N} \sum_{n=1}^N \left(\frac{1}{\ln(\sum_{k=1}^K \exp(\mathbf{w}_k^T \mathbf{x}_n))} \cdot \exp(\mathbf{w}_i^T \mathbf{x}_n) \cdot \mathbf{x}_n - \mathbb{I}_{y_n = i} \mathbf{x}_n \right) \\
&= \frac{1}{N} \sum_{n=1}^N ((h_i(\mathbf{x}_n) - \mathbb{I}_{y_n = i}) \mathbf{x}_n)
\end{aligned}$$

4.



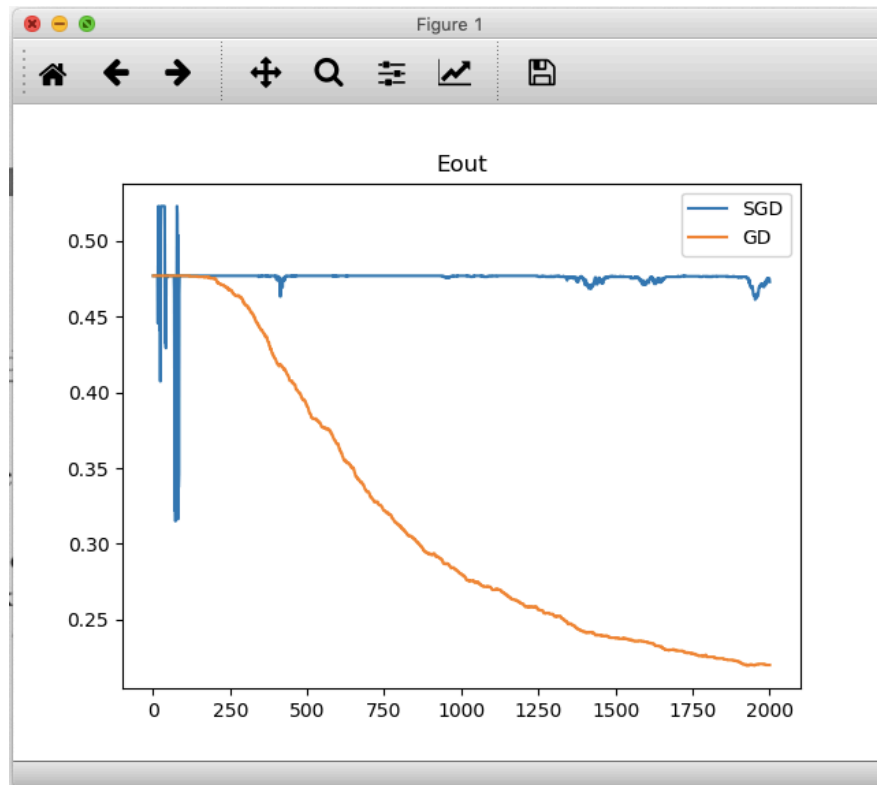
SGD stands for Stochastic Gradient Descent, GD stands for Gradient Descent

We can see that the Ein of GD decrease steadily, which is reasonable because GD calculates the ‘true’ gradient so the Ein of GD should go down steadily. SGD, on the other hand, calculates the gradient on one point to approximate the ‘true’ gradient. This causes the Ein of SGD to not decrease steadily, because sometimes the ‘approximate’ gradient is far away from the true gradient. However, after enough iterations, SGD still can reach the result of GD. In fact, I changed the iterations to 20000 to test my hypothesis. Here’s the result.



For 2500~20000 iterations, the Ein of GD almost doesn’t change, so I think there’s some noise in the data.

5.



SGD stands for Stochastic Gradient Descent, GD stands for Gradient Descent

The Eout graph is similar to the Ein graph, but there are still some differences. After 2000 iterations, Eout of GD is close but slightly greater than Ein of GD. This is reasonable because of VC Bound which tells us that the distance between Ein and Eout will not be too far.

6.

$$\text{Let } \mathbf{h} = \begin{bmatrix} h_1(\mathbf{x}_1) & h_2(\mathbf{x}_1) & \cdots & h_K(\mathbf{x}_1) \\ h_1(\mathbf{x}_2) & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & h_{K-1}(\mathbf{x}_{N-1}) \\ h_1(\mathbf{x}_N) & \cdots & h_{K-1}(\mathbf{x}_N) & h_K(\mathbf{x}_N) \end{bmatrix}, \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix}, \mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_N \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

$$RMSE(H) = \sqrt{\frac{1}{N} \sum_{n=1}^N \left(y_n - \sum_{k=1}^K w_k h_k(\mathbf{x}_n) \right)^2} = \sqrt{\frac{1}{N} \|\mathbf{y} - \mathbf{h}\mathbf{w}\|^2}$$

$$\underset{\mathbf{w}}{\operatorname{argmin}} RMSE(H) = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{h}\mathbf{w}\|^2 = \underset{\mathbf{w}}{\operatorname{argmin}} (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{h}^T \mathbf{y} + \mathbf{w}^T \mathbf{h}^T \mathbf{h} \mathbf{w})$$

Let $err(\mathbf{w}) = \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{h}^T \mathbf{y} + \mathbf{w}^T \mathbf{h}^T \mathbf{h} \mathbf{w}$, the gradient of $err(\mathbf{w})$ is:

$$\nabla err(\mathbf{w}) = 2\mathbf{h}^T \mathbf{h} \mathbf{w} - 2\mathbf{h}^T \mathbf{y}$$

To minimize $err(\mathbf{w})$, $\nabla err(\mathbf{w})$ should be zero, therefore:

$$\mathbf{w} = (\mathbf{h}^T \mathbf{h})^{-1} \mathbf{h}^T \mathbf{y}$$

Now, we need to get $\mathbf{h}^T \mathbf{y}$ to get \mathbf{w} .

How do we get $\mathbf{h}^T \mathbf{y}$?

$$\mathbf{h}^T \mathbf{y} = \begin{bmatrix} h_1(\mathbf{x}_1) & h_1(\mathbf{x}_2) & \cdots & h_1(\mathbf{x}_N) \\ h_2(\mathbf{x}_1) & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & h_{K-1}(\mathbf{x}_N) \\ h_K(\mathbf{x}_1) & \cdots & h_{K-1}(\mathbf{x}_{N-1}) & h_K(\mathbf{x}_N) \end{bmatrix} \mathbf{y}$$

Let \mathbf{r}_i be the i^{th} row of \mathbf{h}^T , and $\mathbf{r}_0 = [h_0(\mathbf{x}_1) \ h_0(\mathbf{x}_2) \ \cdots \ h_0(\mathbf{x}_N)] = [\mathbf{0} \ \mathbf{0} \ \cdots \ \mathbf{0}]$

From the problem statement,

$$e_i = \sqrt{\frac{1}{N} \sum_{n=1}^N (y_n - h_i(\mathbf{x}_n))^2} = \sqrt{\frac{1}{N} \|\mathbf{y} - \mathbf{r}_i^T\|^2}$$

$$Ne_i^2 = \|\mathbf{y} - \mathbf{r}_i^T\|^2 = \mathbf{y}^T \mathbf{y} - 2\mathbf{r}_i \mathbf{y} + \mathbf{r}_i \mathbf{r}_i^T$$

For any $j \in [1, K]$,

$$N(e_0^2 - e_1^2) = 2(\mathbf{r}_1 - \mathbf{r}_0) \mathbf{y} + \mathbf{r}_0 \mathbf{r}_0^T - \mathbf{r}_1 \mathbf{r}_1^T$$

$$N(e_1^2 - e_2^2) = 2(\mathbf{r}_2 - \mathbf{r}_1) \mathbf{y} + \mathbf{r}_1 \mathbf{r}_1^T - \mathbf{r}_2 \mathbf{r}_2^T$$

\vdots

$$N(e_{j-1}^2 - e_j^2) = 2(\mathbf{r}_j - \mathbf{r}_{j-1}) \mathbf{y} + \mathbf{r}_{j-1} \mathbf{r}_{j-1}^T - \mathbf{r}_j \mathbf{r}_j^T$$

Add them all, we get:

$$N(e_0^2 - e_j^2) = 2(\mathbf{r}_j - \mathbf{r}_0) \mathbf{y} + \mathbf{r}_0 \mathbf{r}_0^T - \mathbf{r}_j \mathbf{r}_j^T = 2\mathbf{r}_j \mathbf{y} - 2\mathbf{r}_0 \mathbf{y} + \mathbf{r}_0 \mathbf{r}_0^T - \mathbf{r}_j \mathbf{r}_j^T$$

Notice that $\mathbf{r}_0 \mathbf{y} = 0, \mathbf{r}_0 \mathbf{r}_0^T = 0$, So:

$$Ne_0^2 - Ne_j^2 = 2\mathbf{r}_j \mathbf{y} - \mathbf{r}_j \mathbf{r}_j^T$$

$$\mathbf{r}_j \mathbf{y} = \frac{1}{2} (Ne_0^2 - Ne_j^2 + \mathbf{r}_j \mathbf{r}_j^T) = \frac{1}{2} (Ne_0^2 - Ne_k^2 + \|\mathbf{r}_j\|^2)$$

And then we get:

$$\mathbf{h}^T \mathbf{y} = \begin{bmatrix} h_1(\mathbf{x}_1) & h_1(\mathbf{x}_2) & \cdots & h_1(\mathbf{x}_N) \\ h_2(\mathbf{x}_1) & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & h_{K-1}(\mathbf{x}_N) \\ h_K(\mathbf{x}_1) & \cdots & h_{K-1}(\mathbf{x}_{N-1}) & h_K(\mathbf{x}_N) \end{bmatrix} \mathbf{y} = \begin{bmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \\ \vdots \\ \mathbf{r}_K \end{bmatrix} \mathbf{y} = \sum_{k=1}^K \mathbf{r}_k \mathbf{y}$$

$$\text{By and large, when } \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix} = (\mathbf{h}^T \mathbf{h})^{-1} \mathbf{h}^T \mathbf{y} =$$

$$\left(\begin{bmatrix} h_1(\mathbf{x}_1) & h_1(\mathbf{x}_2) & \cdots & h_1(\mathbf{x}_N) \\ h_2(\mathbf{x}_1) & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & h_{K-1}(\mathbf{x}_N) \\ h_K(\mathbf{x}_1) & \cdots & h_{K-1}(\mathbf{x}_{N-1}) & h_K(\mathbf{x}_N) \end{bmatrix} \begin{bmatrix} h_1(\mathbf{x}_1) & h_2(\mathbf{x}_1) & \cdots & h_K(\mathbf{x}_1) \\ h_1(\mathbf{x}_2) & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & h_{K-1}(\mathbf{x}_{N-1}) \\ h_1(\mathbf{x}_N) & \cdots & h_{K-1}(\mathbf{x}_N) & h_K(\mathbf{x}_N) \end{bmatrix} \right)^{-1} \sum_{k=1}^K \left(\frac{1}{2} (Ne_0^2 - Ne_k^2 + \sum_{n=1}^N h_k(\mathbf{x}_n)^2) \right)$$